# Optimizers
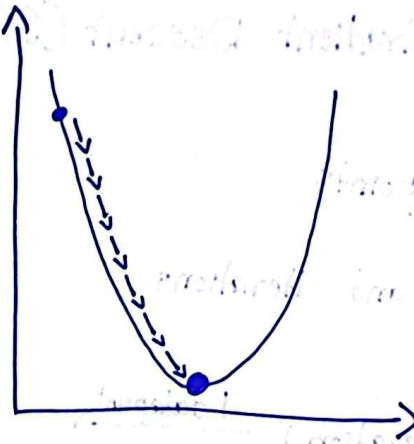
## → Gradient Descent

$$W_{new} = W_{old} - \alpha \frac{\partial L}{\partial W_{old}}$$



### MSE

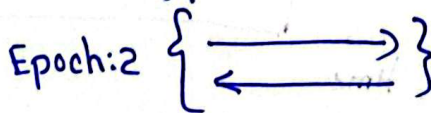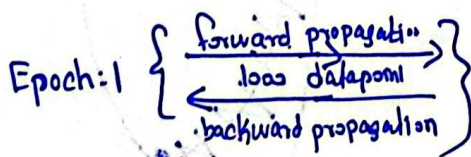Loss function $= (y - \hat{y})^2$

Cost function $= \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$

### Epochs, Iterations

Epoch:1 $\left\{ \begin{array}{c} \text{forward propagation} \\ \text{1000 datapoint} \\ \text{backward propagation} \end{array} \right\}$

Epoch:2 $\left\{ \longrightarrow \atop \longleftarrow \right\}$

Epoch: n $\left\{ \longrightarrow \atop \longleftarrow \right\}$

Dataset = 1000 datapoints
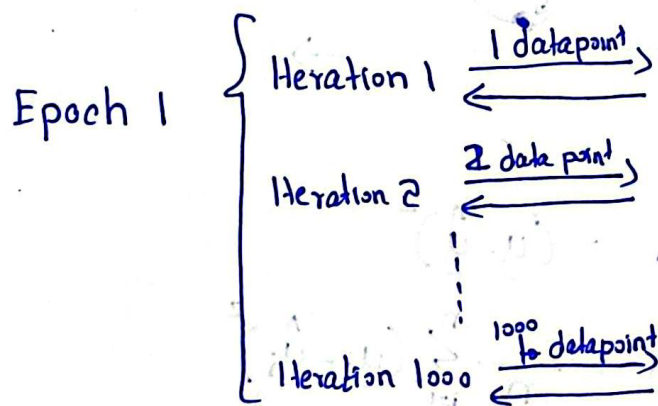
1 Epoch = 1 Iteration

> Advantages

① Convergence will happen

> Disadvantage

① Huge resource required such as RAM and GPU

---> Stochastic Gradient Descent (SGD)
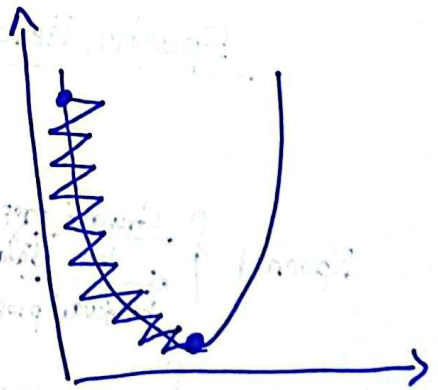
1000 datapoints

Epochs and Iterations

Epoch 1 { Iteration 1 $\xrightarrow{\text{1 datapoint}}$ $\xleftarrow{}$

Iteration 2 $\xrightarrow{\text{2 data point}}$ $\xleftarrow{}$

Iteration 1000 $\xrightarrow{\text{1000}^{\text{th}} \text{ datapoint}}$ $\xleftarrow{}$

> Advantages

① Resource issue solved

> Disadvantages

① Time complexity

② Convergence will take more time

③ Noise gets introduced

# → Mini Batch SGD

Epoch, Iteration, Batch size

Datapoints = 100000

Batch-size = 1000

$$\text{No. of iterations} = \frac{\text{Datapoints}}{\text{batch size}}$$

$$= \frac{100000}{1000}$$

$$= 100$$

Epoch 1 : Iteration : 1

    1000 datapoints

$$\text{Cost function} = \sum_{i=1}^{1000} (y_i - \hat{y}_i)^2$$
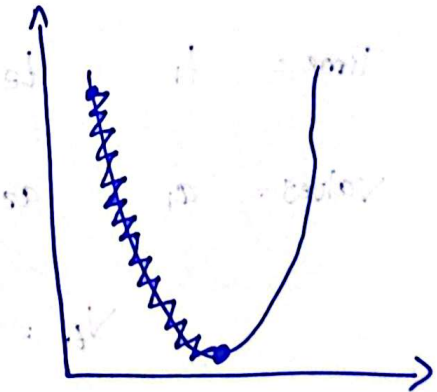
Epoch 1 : Iteration : 2

    1000 datapoints

    ⋮

Epoch 1 : Iteration : 100

    1000 datapoints

## > Advantages

① Converges faster than SGD

② Noise will be less than SGD

③ Efficient resource usage

## > Disadvantage

① Noise still exists

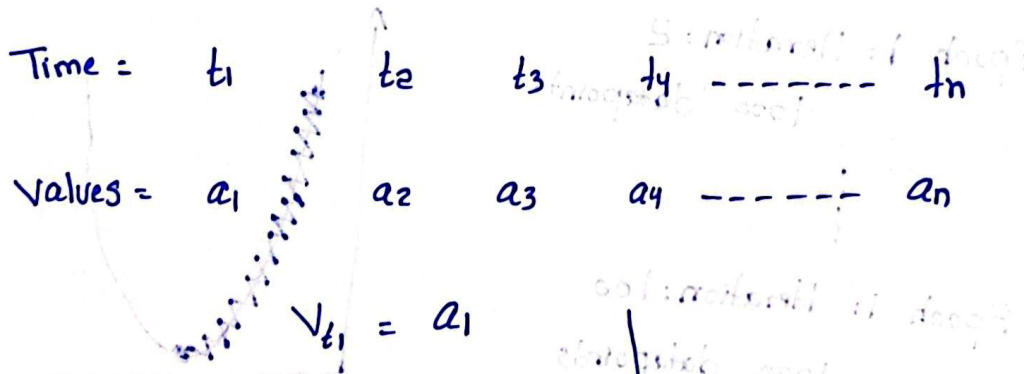## $\longrightarrow$ SGD with Momentum

$$W_{new} = W_{old} - \alpha \cdot \frac{\partial L}{\partial W_{old}}$$

$$b_{new} = b_{old} - \alpha \frac{\partial L}{\partial b_{old}}$$

$$W_t = W_{t-1} - \alpha \frac{\partial L}{\partial W_{t-1}}$$
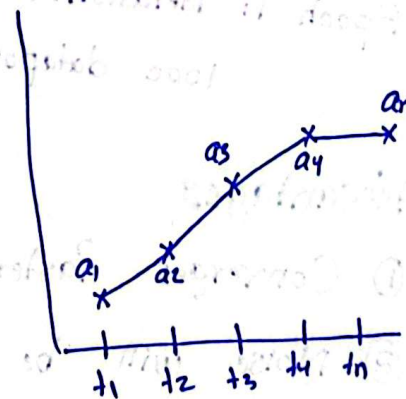
### Exponential Weight Average

Use to perform smoothening

Time = $t_1$     $t_2$     $t_3$    $t_4$ ------- $t_n$

Values = $a_1$     $a_2$    $a_3$    $a_4$ ------ $a_n$

$$V_{t_1} = a_1$$

$$V_{t_2} = \beta V_{t_1} + (1-\beta) a_2$$

$$\beta = 0.95$$
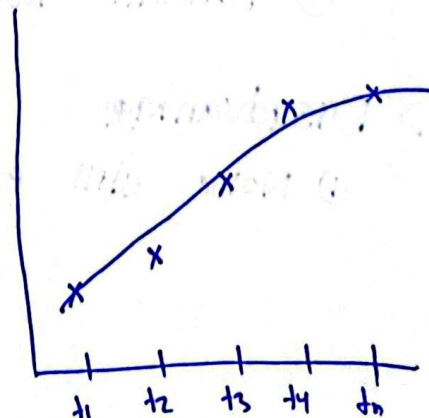
$$V_{t_2} = 0.95 a_1 + (1-0.95) a_2$$

$$= 0.95 a_1 + 0.05 a_2$$

> if $\beta$ has high value,
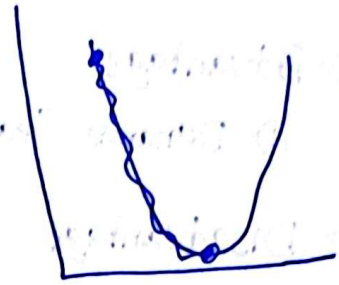> it will control the previous
> value more.

$$V_{t_3} = \beta V_{t_2} + (1-\beta) a_3$$

$$= 0.95 \left[ 0.95 a_1 + 0.05 a_2 \right] + 0.05 a_3$$
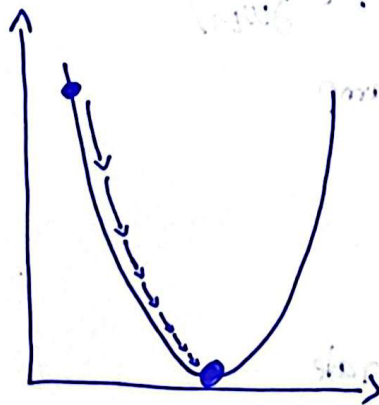
> Advantage
  ① Reduce the noise
  ② Quick convergence



——> Adagrad : Adaptive Gradient Descent

$$W_{new} = W_{old} - \alpha \frac{\partial L}{\partial W_{old}}$$

$$\alpha = fixed$$

$$\Downarrow$$

$$\alpha = dynamic$$



As the convergences happen,
the learning rate will
change

$$W_{new_t} = W_{old_{t-1}} - \alpha' \frac{\partial L}{\partial W_{old_{t-1}}}$$

$$\alpha' = \frac{\alpha}{\sqrt{\alpha_{t_{new}} + \epsilon}} \longrightarrow epsilon$$

$$\alpha_t = \sum_{i=1}^{t} \left( \frac{\partial L}{\partial W_t} \right)^2 \qquad \therefore \alpha_t \text{ increases, } \alpha' \text{ decreases}$$

$$t=1 \qquad\qquad t=2 \qquad\qquad t=3$$
$$\alpha = 0.01 \qquad \alpha = 0.005 \qquad \alpha = 0.003$$

> Advantage:

① Dynamic learning rate

> Disadvantage

① Possibility of learning rate to become approx. zero

② Convergence may never occur

---→ Adadelta and RMSPROP

$$\alpha' = \frac{\alpha}{\sqrt{Sdw + \epsilon}}$$

$$Sdw_t = \beta \, Sdw_{t-1} + (1-\beta)\left(\frac{\partial L}{\partial w_{t-1}}\right)^2$$

for the first time stamp

$$Sdw_t = 0$$

> Advantages

① Dynamic learning rate

② Smoothening Exponential Weighted Average

$$w_t = w_{t-1} - \alpha' \frac{\partial L}{\partial w_{t-1}}$$

$\longrightarrow$ **Adam Optimizer**

SGD with Momentum + RMSPROP

$$W_t = W_{t-1} - \alpha' \, Vdw \qquad \Rightarrow \text{weight updation}$$

$$b_t = b_{t-1} - \alpha' \, Vdb \qquad \Rightarrow \text{bias updation}$$

$$\alpha' = \frac{\alpha}{\sqrt{Sdw} + \epsilon}$$

$$Sdw_t = 0$$

$$Sdw_t = \beta \, Sdw_{t-1} + (1-\beta)\left(\frac{\partial L}{\partial w_{t-1}}\right)^2$$

$$\left. \begin{array}{l} Vdw_t = \beta \, Vdw_{t-1} + (1-\beta)\dfrac{\partial L}{\partial w_{t-1}} \\[3mm] Vdb_t = \beta \, Vdb_{t-1} + (1-\beta)\dfrac{\partial L}{\partial b_{t-1}} \end{array} \right\} \begin{array}{l} \text{Momentum} \\ \quad\hookrightarrow \text{Smoothening} \end{array}$$